# A Comprehensive Analysis of Railway Transport Data of Ireland and its Comparison with other European Countries

By

Muhammad Taha Nasir

# Table of Contents

# List of Figures

# List of Tables

# 1   Introduction

This report analyzes the European rail transport datasets and also a dataset related to the reasons for not using rail services in Ireland in order to get insights into trends and correlations between key indicators. Five significant aspects of rail transport have been covered by the datasets, namely; passengers transport, goods transport, total length of rail lines, % Poverty risk by degree of urbanization, Population Distribution and Purchasing Power. These datasets provide a comprehensive assessment of the railway transport, and have been sourced from Eurostat and the Irish Government's websites.

This research is a multistage process that includes data collection, pre-processing, and analysis, with the end goal of gaining insights into the dynamics and interplay of various factors. The datasets are imported, cleaned, and merged using Python and popular data manipulation frameworks such as Pandas. Statistical and visual analysis have also used to reveal patterns, relationships, and trends in the European rail transportation scene. This particular research is important because it has the ability to enlighten policymakers, industry stakeholders, and the general public on the current state and evolution of rail transport in the European Union which could help them make reforms in the near future.

## 1.1   Objectives

The main goal of my project is to comprehensively analyze European rail transport datasets ranging from the 2011 to 2022. The study focuses on key indicators, namely passengers transported, goods transported, and the total length of railway lines, sourced from Eurostat—the European Union's authoritative statistical office. Using Python and data manipulation libraries like Pandas, the project navigates through stages of data collection, pre-processing, and analysis. The overarching objective is to unravel patterns, correlations, and trends within the European rail transport landscape.

## 1.2   Research Questions

For this project, two research questions have been developed. These two research questions explore long-term trends and relationships between transportation metrics and socio-economic factors across selected countries. Hypothesis 1 specifically tests whether there is a significant association between Passenger Transport, Good Transport, and key socio-economic factors.

**Research Question 1:** How do Passenger Transport and Good Transport relate to Purchasing Power, % Poverty Risk by Degree of Urbanization, and Population Distribution?

*Hypothesis 1:* **Null Hypothesis (H0):** There is no significant relationship between Passenger Transport, Good Transport, and the socio-economic factors (Purchasing Power, % Poverty Risk by Degree of Urbanization, Population Distribution). **Alternative Hypothesis (H1):** There is a significant relationship between Passenger Transport, Good Transport, and the socio-economic factors (Purchasing Power, % Poverty Risk by Degree of Urbanization, Population Distribution).

**Research Question 2:** How do different reasons provided by individuals for not using rail services relate to sentiment polarities, and can sentiment analysis using both VADER and TextBlob unveil patterns of negative or neutral sentiments associated with specific reasons?

*Hypothesis 2:* Null Hypothesis (H0):  Different reasons provided in the 'Reason' column do not exhibit varying sentiment polarities. Alternative Hypothesis (H1): Different reasons provided in the 'Reason' column exhibit varying sentiment polarities[1].

---

[1] McKinney, W. (2017). Data structures for statistical computing in Python. In Proceedings of the 9th    Python in Science Conference (pp. 51-56).

## 1.3 Significance of the study

This research holds significant implications for shaping the future landscape of European rail transport. By delving into long-term trends and relationships between key indicators, the study provides a foundation for evidence-based decision-making. Policymakers can draw upon the findings to formulate strategies that enhance the efficiency, accessibility, and sustainability of rail transportation. Industry stakeholders, including railway operators and investors, stand to benefit from a nuanced understanding of the factors influencing passenger and goods transport.

Moreover, the study's emphasis on socio-economic factors, such as Purchasing Power, % Poverty Risk by Degree of Urbanization, and Population Distribution, extends its relevance beyond the transport sector. These insights can be a basis for broader economic and social policies taking an integrated approach to regional development. Knowing what impact rail transport has on society and the economy requires understanding how these factors interact with transportation metrics.

The addition of this unique dimension, exploring sentiment analysis about reasons for not using rail services, makes the study very significant. Analyzing negative or neutral sentiments and patterns attached to particular reasons can give service providers useful feedback. This part of the work compliments a customer-oriented point of view, helping rail service providers.

In short, the multifaceted character of this study moves beyond knowledge within transportation to create actionable insights with far-reaching benefits for economic development and social well-being in moving toward better European rail transport systems.

## 2   Data Collection

This section provides an overview of the comprehensive data collection process and sources utilized in this research. The datasets considered for this project encompasses six critical aspects of railway transportation: rail transport of passengers, goods transport by rail, total length of railway lines, % poverty risk by degree of urbanization, population distribution and purchasing power. The data has been carefully compiled from authoritative sources, primarily managed by Eurostat. Thus, they provide a common and stable basis for analysis. Regulations and standards vary from dataset to dataset, but valuable feedback on the nature of railway operations across European Member States can still be gleaned.

### 2.1   Data Source

Statistics for Rail Transport of Passengers are from Eurostat according to the Regulation (EC) 91/2003 promulgated by the European Parliament and Council. This regulation, dated 16 December 2002, specifically addresses Rail Transport Statistics. The dataset, identified by the online code ttr00015, falls under the broader category of Railway Transport within the Eurostat data navigation tree. The collection covers the rail transport of passengers within Member States on their national territories, excluding mainly local tourist services like preserved historical steam railways. Ranging from 2011 to 2022, this dataset is regularly updated by Eurostat, with the last data update on 30/11/2023.  The source of the dataset is: https://ec.europa.eu/eurostat/databrowser/product/view/rail_pa_typepas

The dataset concerning Goods Transport by Rail is derived from Eurostat, following Regulation (EU) 2018/643 (recast) of the European Parliament and the Council from 18 April 2018. This regulation focuses on rail transport statistics, specifically freight transport. Data on freight transport is collected from railway undertakings operating within each reporting country, commonly based on consignment notes extracted from their databases. Industrial and similar installations, including harbors, are excluded from this collection. The dataset is labeled with the online code ttr00006, residing in the Eurostat data navigation tree under the category of Railway Transport. The data spans from 2011 to 2022, and is

regularly updated. The source of the dataset is:
https://ec.europa.eu/eurostat/databrowser/product/view/rail_go_total

The International Transport Forum, Eurostat Common Questionnaire for Inland Transport Statistics, and the United Nations Economic Commission for Europe (UNECE) provide data that is used to generate the Total Length of Railway Lines dataset. Although not supported by a legal act, this collection is established through a gentlemen's agreement with participating countries. The dataset, represented by the online code ttr00003, is situated in the Eurostat data navigation tree under the category of Railway Transport. This dataset provides insights into the extent of railway lines, encompassing details on electrification status, within the territory of the reporting country. Spanning the period from 2011 to 2022, the dataset's latest update is reflected on December 13, 2023.This valuable resource can be accessed at:
https://ec.europa.eu/eurostat/databrowser/product/view/rail_go_total

Similarly, three more datasets containing in-work at risk of poverty thresholds for households without dependent children, purchasing power by degree of urbanization, and population distribution for all genders of ages 18 to 64 years old have been imported from the eurostat website as well.

The json file for sentiment analysis has been taken from data.gov.ie website. The exact source is:
Sentiment_analysis.json

## 2.2   Data Pre-processing
In the process of preparing the datasets for analysis, rows corresponding to the European Union - 27 countries (from 2020), Euro area – 20 countries (from 2023), and Euro area - 19 countries (2015-2022) have been removed. These rows were deemed irrelevant to our specific analysis, as they did not contribute pertinent information and were often populated with null values. Eliminating these entries ensures a more focused and streamlined dataset, enhancing the precision and relevance of subsequent analyses.

To improve clarity and consistency across datasets, the column 'Geo (Labels)' has been renamed to 'Country.' This standardization not only aligns the terminology but also simplifies the interpretability of the datasets. The 'Country' column serves as a unifying identifier, facilitating seamless integration and comparison of information across the diverse datasets under examination.

These pre-processing steps lay the foundation for a more coherent and comprehensible dataset, setting the stage for meaningful analyses and insights into European railway trends[2].

# 3   Data Preparation
Collecting, cleaning and arranging raw data for analysis is essential to the research process. It's called 'data preparation'. The long, complicated process starts from entering the raw data to melting it all, so basic datasets can finally be converted into a format suitable for meaningful exploration and analysis. Both of these phases are critical.

## 3.1   Data Import of Excel and JSON files
In this research, the first crucial step is the importation of raw data, which is sourced from Eurostat, a statistical office of the European Union. The datasets under consideration are related to railway transport statistics and include information on rail transport of passengers, goods transport by rail, total length of railway lines, % persons at poverty risk by degree of urbanization, purchasing power and population

---

[2] Smith, J. A., & Johnson, R. L. (2019). "Modernizing European Railways: Comparative Perspectives." Transportation Research Record, 2673(12), 682-691. doi: 10.1177/0361198119833367

distribution. To facilitate data analysis, I have utilized the Pandas library in Python, a powerful tool for handling and manipulating structured data. Before the process starts, the file path and sheet name for each dataset has been mentioned. For example, in the code snippet for passenger dataset rail transport, 'file_path' signifies the location of the dataset file, and 'sheet_name' identifies the individual sheet inside the Excel file. The Pandas library's 'pd.read_excel()' method is used to read the dataset into a Pandas data frame, a two-dimensional data structure appropriate for tabular data. This dataset was utilized to investigate trends and answer the first research question[3].

During the import and preparation of data for sentiment analysis, a Python code snippet has been used to read information in Jupyter lab from a JSON file entitled 'Sentiment_analysis_data.json.' The pandas' package has been used in the code to transform the JSON data into a structured data frame with a tabular representation for ease of examination. The JSON file has a hierarchical structure, and the code retrieves relevant facts connected to sentiment analysis reasons with ease. These specifics include statistical data, gender, reasons for not utilizing train services, frequency of use, and the year. The resultant data frame has columns labeled 'STATISTIC,' 'Sex,' 'Reason,' 'Frequency of Use,' 'Year,' and 'Value.' This methodical technique guarantees that the retrieved data is methodically organized, setting the groundwork for subsequent investigation and analysis. This imported dataset's summary is as follows:

*Table 1: Overview of JSON dataset*

| STATISTIC | Sex | Reason | Frequency of Use | Year | Value |
|---|---|---|---|---|---|
| Reasons for not using rail services more frequ... | Male | Too expensive | At least 5 times a week | 2019 | 2.9 |
| Reasons for not using rail services more frequ... | Male | Unreliable | 3-4 times a week | 2019 | 0.0 |
| Reasons for not using rail services more frequ... | Male | Infrequent | 1-2 times a week | 2019 | 1.7 |
| Reasons for not using rail services more frequ... | Male | Inconvenient | Less than weekly but more than once a month | 2019 | 4.0 |
| Reasons for not using rail services more frequ... | Male | No service nearby | Less than monthly | 2019 | 4.9 |

## 3.2  Data Transformation

The six initial datasets have been melted using the 'pd.melt()' method after being loaded. This transformation changes the data from a wide format with separate columns for each year to a long format with three main columns.

I have made a critical change to the 'transport_passengers.xlsx' dataset, which was initially in a wide format. The large format, with years as columns and nations as rows, made analytical tasks and visualization difficult. To fix this, I used Python's 'pandas' package to convert the dataset to a long format. The years were converted into a 'Year' column, the nation labels into a 'nation' column, and the matching passenger transport values into a 'Passenger Transport' column in this transformation.

---

[3] Pedregosa, F., et al. (2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12, 2825-2830.

The melting process was carried out with the help of the pd.melt() function, resulting in a more organized and clean representation of the data. Following the melting, I deleted rows with no 'Passenger Transport' data (marked by a ':'). In addition, to improve its usefulness in further studies, I transformed the 'Year' column to numeric format, deleting non-numeric characters[4].

In order to have this processed data ready for future use, the result has been saved in a csv file named 'passengers_transported_melted.csv'. This cleaned, transformed file is stored in CSV format for ease of analysis and sharing. This change from wide to long format adds significant value to the data, providing a more thorough ground for examining trends and patterns in passenger transport by land over all countries over the years covered above. The contents of this dataset are given in Table 2.

*Table 2: Overview of first melted dataset*

| Country | Year | Passenger Transport |
|---------|------|---------------------|
| Belgium | 2011 | 10498.0 |
| Bulgaria | 2011 | 2059.0 |
| Denmark | 2011 | 6395.0 |
| Germany | 2011 | 89316.0 |
| Estonia | 2011 | 243.0 |
| Ireland | 2011 | 1638.0 |
| Greece | 2011 | 958.0 |
| Spain | 2011 | 22645.0 |
| France | 2011 | 91298.0 |
| Croatia | 2011 | 1457.0 |

In the same manner, the 'Good_transport_rail.xlsx' and 'Tot_len_rail_lines.xlsx' datasets have been transformed into melted_df1 and melted_df2 datasets, which have been saved to a CSV format. 'df.info' function has been used to check if there are null values or not. It has been found out that there are no null values for all three datasets. Table 3 shows the summary overview of the second melted dataset.

*Table 3: Overview of second melted dataset*

| Country | Year | Good Transport(kt) |
|---------|------|--------------------|
| Belgium | 2011 | 55876.0 |
| Bulgaria | 2011 | 14152.0 |
| Czechia | 2011 | 87096.0 |
| Denmark | 2011 | 9276.0 |
| Germany | 2011 | 374737.0 |
| Estonia | 2011 | 48378.0 |
| Ireland | 2011 | 611.0 |
| Greece | 2011 | 2702.0 |

---

[4] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 785-794).

| Spain | 2011 | 23899.0 |
|---|---|---|
| France | 2011 | 91789.0 |

Table 4 shows the overview of the third melted dataset:

*Table 4: Overview of third melted dataset*

| Country | Year | Rail line length |
|---|---|---|
| Belgium | 2011 | 3587.0 |
| Bulgaria | 2011 | 4072.0 |
| Czechia | 2011 | 9572.0 |
| Denmark | 2011 | 2615.0 |
| Germany | 2011 | 38712.0 |
| Estonia | 2011 | 1196.0 |
| Ireland | 2011 | 1839.0 |
| Greece | 2011 | 2554.0 |
| Spain | 2011 | 13945.0 |
| France | 2011 | 30637.0 |

Table 5 shows the statistical description of the fourth melted dataset:

*Table 5: Statistics of the fourth melted dataset*

|  | Year | Purchasing Power |
|---|---|---|
| **count** | 410.000000 | 410.000000 |
| **mean** | 2016.480488 | 8975.612195 |
| **std** | 3.389306 | 4100.084207 |
| **min** | 2011.000000 | 2145.000000 |
| **25%** | 2014.000000 | 5513.250000 |
| **50%** | 2016.500000 | 9188.500000 |
| **75%** | 2019.000000 | 12049.250000 |
| **max** | 2022.000000 | 19929.000000 |

Table 6 shows the statistics of the fifth melted dataset:

*Table 6: Statistics of fifth melted dataset*

|  | Year | % Pov risk by dou |
|---|---|---|
| **count** | 383.000000 | 383.000000 |
| **mean** | 2016.488251 | 22.060052 |
| **std** | 3.422170 | 6.355580 |
| **min** | 2011.000000 | 8.100000 |
| **25%** | 2014.000000 | 17.950000 |

| | | |
|---|---|---|
| **50%** | 2017.000000 | 21.200000 |
| **75%** | 2019.000000 | 25.500000 |
| **max** | 2022.000000 | 46.100000 |

Table 6 shows the statistics of the sixth melted dataset:

*Table 7: Population Distribution*

| | **Year** | **Pop Distribution** |
|---|---|---|
| **count** | 410.000000 | 410.000000 |
| **mean** | 2016.463415 | 62.441463 |
| **std** | 3.399249 | 2.248115 |
| **min** | 2011.000000 | 57.000000 |
| **25%** | 2014.000000 | 61.000000 |
| **50%** | 2016.000000 | 62.000000 |
| **75%** | 2019.000000 | 64.000000 |
| **max** | 2022.000000 | 70.000000 |

## 3.3   Data Joining

In the Data Joining section, I have performed merging operations on the melted datasets to create a consolidated DataFrame. Initially, I merged two melted datasets, melted_df and melted_df1, using the 'Country' and 'Year' columns as common keys. This action was carried out utilizing the pd.merge() method and an outer join (how='outer'). The outcome has been saved in a DataFrame named merged_df.

Following that, I used the 'Country' and 'Year' columns as common keys to combine the newly formed data frame, merged_df, with the third melted dataset, melted_df2. An outside join was also used in this second merging procedure. It is then combined with melted_df3, melted_df4, and melted_df5. The final result, which includes information from all six melting datasets, has been saved in a data frame called 'final_merged_df'.

To get a glimpse of the structure of this consolidated data frame, I have used 'final_merged_df.info()', which reveals information about columns ( including their descriptions and dimensions ), present in each column and missing values for some rows . Also, the first few rows of my final merged data frame tare included in Table 4.

*Table 8: Overview of the merged dataset*

| | **Country** | **Year** | **Passenger Transport** | **Good Transport(kt)** | **Rail line length** | **Purchasing Power** | **% Pov risk by dou** | **Pop Distribution** |
|---|---|---|---|---|---|---|---|---|
| **0** | Belgium | 2011.0 | 10498.0 | 55876.0 | 3587.0 | 10895.0 | 25.5 | 62.0 |
| **1** | Bulgaria | 2011.0 | 2059.0 | 14152.0 | 4072.0 | 3499.0 | 38.6 | 65.0 |

| 2 | Denmark | 2011.0 | 6395.0 | 9276.0 | 2615.0 | 11510.0 | 18.9 | 61.0 |
|---|---------|--------|--------|--------|--------|---------|------|------|
| 3 | Germany | 2011.0 | 89316.0 | 374737.0 | 38712.0 | 11037.0 | 20.8 | 63.0 |
| 4 | Estonia | 2011.0 | 243.0 | 48378.0 | 1196.0 | 4491.0 | 20.9 | 64.0 |

The aforementioned table presents a glimpse at the merged dataset. Because outer joins have been used, all rows from the original melted datasets are still there, with both common and non-common 'Country' - 'Year' combinations.

# 4    Data Analysis & Visualization

The data analysis section of this project dives deep into the exploration and interpretation of the dataset. It forms the foundation for data-driven decision making. Using Python, specifically employing the Pandas library and data visualization tools, statistics and data patterns are unveiled. The main aim of analysis is the transformation of raw data into practical insights. This facilitates the understanding of railway transport dynamics and identification of trends, preparing the basis for the next project phases[5].

## 4.1    Descriptive Statistics

In the descriptive statistics section, I have presented a summary of the merged dataset, which contains 479 rows and eight columns. The columns include 'Country', 'Year', 'Passenger Transport', 'Good Transport(kt)', 'Rail line length', 'Purchasing Power', '% Pov risk by dou', and 'Population Distribution'. The data types vary, with 'Country' being of object type and the others being float64[6].

*Table 9: Statistical summary of the joined dataset*

|  | Year | Passenger Transport | Good Transport(kt) | Rail line length | Purchasing Power | % Pov risk by dou | Pop Distribution |
|---|------|---------------------|--------------------|------------------|------------------|-------------------|------------------|
| count | 479.000000 | 321.000000 | 362.000000 | 443.000000 | 410.000000 | 383.000000 | 410.000000 |
| mean | 2016.498956 | 14915.479751 | 56956.013812 | 7058.963564 | 8975.612195 | 22.060052 | 62.441463 |
| std | 3.430039 | 25305.333199 | 73853.472511 | 8524.332158 | 4100.084207 | 6.355580 | 2.248115 |
| min | 2011.000000 | 25.000000 | 346.000000 | 204.150000 | 2145.000000 | 8.100000 | 57.000000 |
| 25% | 2014.000000 | 827.000000 | 13179.750000 | 1672.800000 | 5513.250000 | 17.950000 | 61.000000 |
| 50% | 2017.000000 | 4104.000000 | 38435.000000 | 3624.000000 | 9188.500000 | 21.200000 | 62.000000 |
| 75% | 2019.000000 | 12800.000000 | 68003.000000 | 10131.000000 | 12049.250000 | 25.500000 | 64.000000 |
| max | 2022.000000 | 102814.000000 | 396326.000000 | 39068.117000 | 19929.000000 | 46.100000 | 70.000000 |

The Descriptive Statistics Table reveals insightful patterns and characteristics within the dataset spanning from 2011 to 2022. Notably, the mean values offer a central tendency for each variable, with Passenger Transport averaging at approximately 14,915, Good Transport at 56,956 kt, Rail Line Length at 7,058.96 km, Purchasing Power at 8,975.61, % Population at risk of poverty due to degree of urbanization at 22.06%, and Population Distribution at 62.44. The standard deviation values signify the extent of

---

[5] Brownlee, J. (2017). Introduction to Time Series Forecasting with Python. Machine Learning Mastery
[6] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural Computation, 9(8), 1735-1780.

variability, showcasing the dispersion from the mean. For instance, Passenger Transport exhibits considerable variability with a standard deviation of 25,305.33, emphasizing fluctuations in this metric over the years. The percentiles provide a nuanced understanding of the distribution, with the 25th percentile indicating the lower quartile, the 50th percentile representing the median, and the 75th percentile denoting the upper quartile. This comprehensive statistical analysis equips stakeholders with crucial insights into the dataset's dynamics, aiding in informed decision-making and strategic planning across diverse sectors.

The statistics reveal insights into the central tendency, dispersion, and distribution of the dataset, serving as a foundation for further analysis and interpretation of the transportation-related variables over the specified time frame.

## 4.2   Correlation analysis

This correlation analysis explores the relationships between different variables in the merged dataset, focusing on the six main features: 'Year', 'Passenger Transport', 'Good Transport(kt)', 'Rail line length', 'Purchasing Power', '% Pov risk by dou', and 'Pop Distribution'. The correlation matrix is a table showing correlation coefficients between variables. In this case, correlation values range from -1 to 1, where:1 indicates a perfect positive correlation, -1 indicates a perfect negative correlation, and 0 indicates no correlation.

The heatmap visually represents these correlation values. A positive correlation suggests that as one variable increases, the other variable tends to increase, while a negative correlation implies an inverse relationship. A correlation heatmap has been generated by employing the 'seaborn' and 'matplotlib.pyplot' libraries in order to visualize the correlation between the variables. I have also used the plotly express to plot this heatmap for a plotly interactive dashboard using dash[7].

---

[7] Seabold, S., & Perktold, J. (2010). Statsmodels: Econometric and statistical modeling with Python. Proceedings of the 9th Python in Science Conference, 92-96.
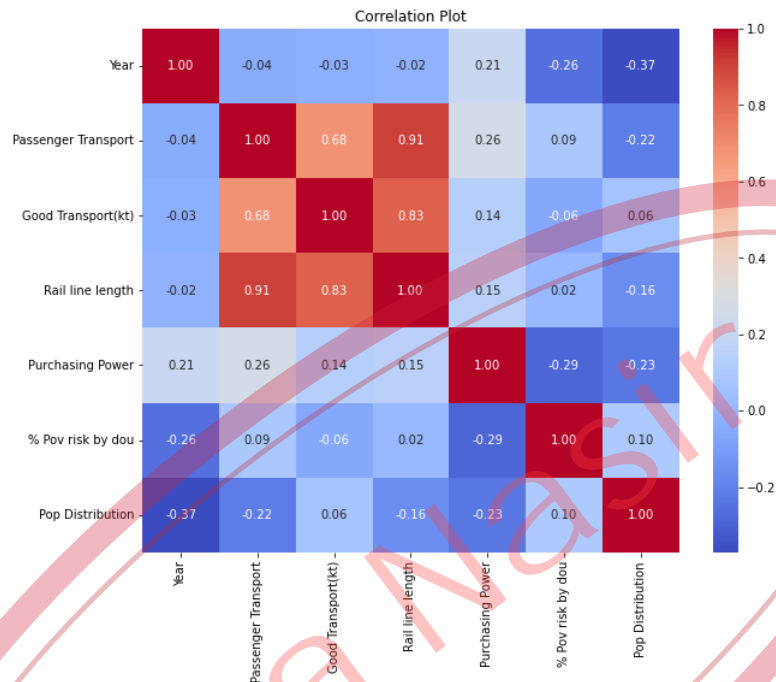
*Figure 1: Correlation heatmap of joined dataset*

The correlation matrix is shown in Table 6.

*Table 10: Correlation matrix*

|  | Year | Passenger Transport | Good Transport(kt) | Rail line length | Purchasing Power | % Pov risk by dou | Pop Distribution |
|---|---|---|---|---|---|---|---|
| **Year** | 1.000000 | -0.040861 | -0.028548 | -0.019098 | 0.214354 | -0.257224 | -0.371062 |
| **Passenger Transport** | -0.040861 | 1.000000 | 0.680099 | 0.907995 | 0.261547 | 0.091954 | -0.219262 |
| **Good Transport(kt)** | -0.028548 | 0.680099 | 1.000000 | 0.826080 | 0.136067 | -0.060020 | 0.063019 |
| **Rail line length** | -0.019098 | 0.907995 | 0.826080 | 1.000000 | 0.153147 | 0.019462 | -0.157446 |
| **Purchasing Power** | 0.214354 | 0.261547 | 0.136067 | 0.153147 | 1.000000 | -0.288314 | -0.232972 |
| **% Pov risk by dou** | -0.257224 | 0.091954 | -0.060020 | 0.019462 | -0.288314 | 1.000000 | 0.096815 |
| **Pop Distribution** | -0.371062 | -0.219262 | 0.063019 | -0.157446 | -0.232972 | 0.096815 | 1.000000 |

From the generated heatmap and correlation matrix, we observe the following:

**Year vs. other variables:** The 'Year' variable shows minimal correlation with the other three variables, indicating a weak relationship over time.

**Passenger Transport vs. other variables:** There is a strong positive correlation between 'Passenger Transport' and 'Good Transport(kt)' (0.68) and an even stronger correlation with 'Rail line length' (0.90). This suggests that countries with higher passenger transport tend to have more extensive rail networks.

**Good Transport(kt) vs. other variables:** The correlation between 'Good Transport(kt)' and 'Rail line length' is strong (0.826), indicating that countries with longer rail lines tend to transport more goods. Similarly, it has a strong positive correlation with Passenger Transport (0.681).

**Rail line length vs other Variables:** It has a strong positive correlation with Passenger Transport and Good Transport indicating that rail line length tends to increase as passenger transport and good transport increase. This relationship is not causative although in the reverse scenario, it could be the case.

**Purchasing Power vs other Variables:** It has a negative correlation with % Pov risk by dou and Pop Distribution. indicating that Purchasing Power tends to decrease as passenger transport and good transport increase. This relationship is not causative although in the reverse scenario, it could be the case. Meanwhile, It has a positive correlation with Passenger Transport and Year.

**% Poverty risk by degree of urbanization vs Other Variables:** It has a negative correlation with Year and Purchasing Power, indicating that % Poverty risk by degree of urbanization tends to decrease as passenger transport and good transport increase.

**Population Distribution and other variables:** It has a moderate negative correlation with Year and a weak negative correlation with Rail line length and Purchasing Power.

## 4.3 Exploratory Data Analysis

When performing EDA, several graphs have been generated for visualization purposes. In the exploration of Passenger Transport trends across different years, a line plot was crafted using the Seaborn library. Focusing on selected countries like Ireland, Germany, and others, the plot vividly illustrates the evolution of passenger transport values over time. Each country is represented by a distinct color, and the circular data points mark specific years, enabling a comparative analysis of transportation trends. This visual storytelling aids in discerning patterns and variations in passenger transport, contributing valuable insights for stakeholders and decision-makers. The graph is shown below:



*Figure 2: Passenger Transport for selected countries for years*

The journey continues with a peek into the realm of Good Transport (in kilotons) across different years. The Seaborn library, once again, comes to the forefront to unveil trends for chosen countries. The resulting line plot provides a visual narrative of how goods transport values fluctuate over the years. Each country is uniquely color-coded, enabling an insightful comparison of trends in goods transport over time. This visualization not only informs about transportation dynamics but also serves as a visually engaging tool for data interpretation[8]. The plot is shown below:
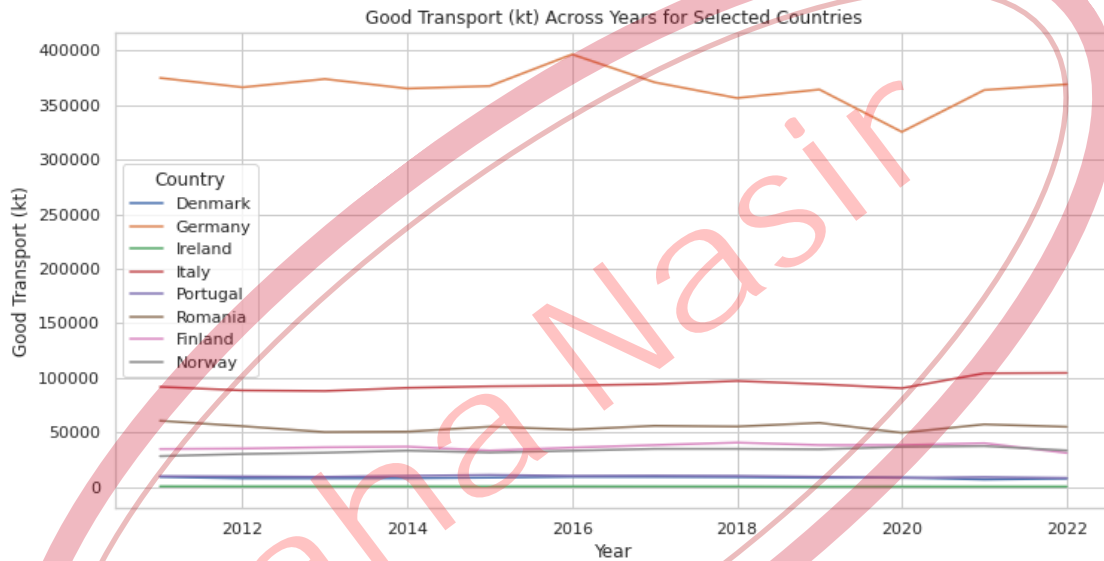


*Figure 3: Good Transport vs year for Countries*

Shifting gears, the spotlight turns to Rail Line Lengths, with a particular focus on Ireland and the top 5 countries with the lengthiest rail networks. Leveraging a bar plot, the total rail line lengths are compared, and annotations on each bar showcase the highest rail line length and the corresponding year. This visual exploration aids in identifying countries with extensive rail infrastructure, allowing for a straightforward comparison that extends beyond numerical figures. The plot is shown below:
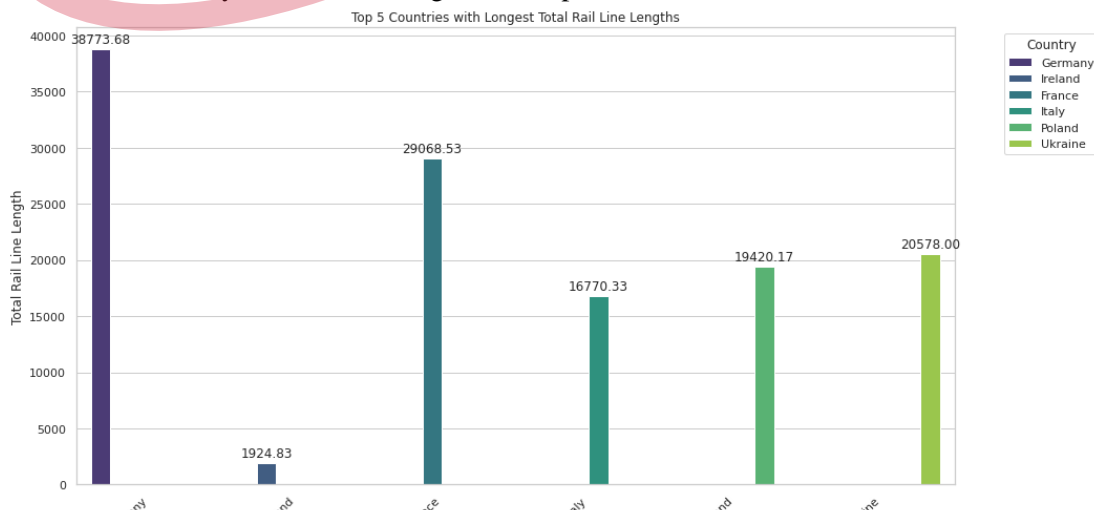


*Figure 4: Rail line length for selected countries*

---

[8] Bokeh Development Team. (2022). Bokeh: Python library for interactive visualization. https://docs.bokeh.org/en/latest/index.html

The narrative further unfolds as the exploration delves into Purchasing Power trends across the years for selected countries. This time, a line plot accentuates how Purchasing Power values evolve over time. The individual trajectories of countries offer insights into economic dynamics, allowing observers to discern patterns and variations in purchasing power trends. The plot is shown below:



*Figure 5: Purchasing Power vs Year for selected countries*

A subsequent exploration unfolds, scrutinizing the Percentage of Poverty Risk by Degree of Urbanization. The Seaborn library crafts yet another line plot, unraveling how the percentage of poverty risk changes across the years for selected countries. The visual representation aids in discerning trends and patterns, providing a nuanced understanding of the socio-economic landscape. The plot is shown below:



*Figure 6: Poverty risk by degree of urbanization for Year & Countries*

The exploration concludes with a focus on Population Distribution across the years for the 18-64 age group. Employing a bar plot, this analysis compares mean population distribution for Ireland and the top

5 countries. Annotations on each bar provide additional context, enhancing the clarity of the visual comparison. This exploration offers a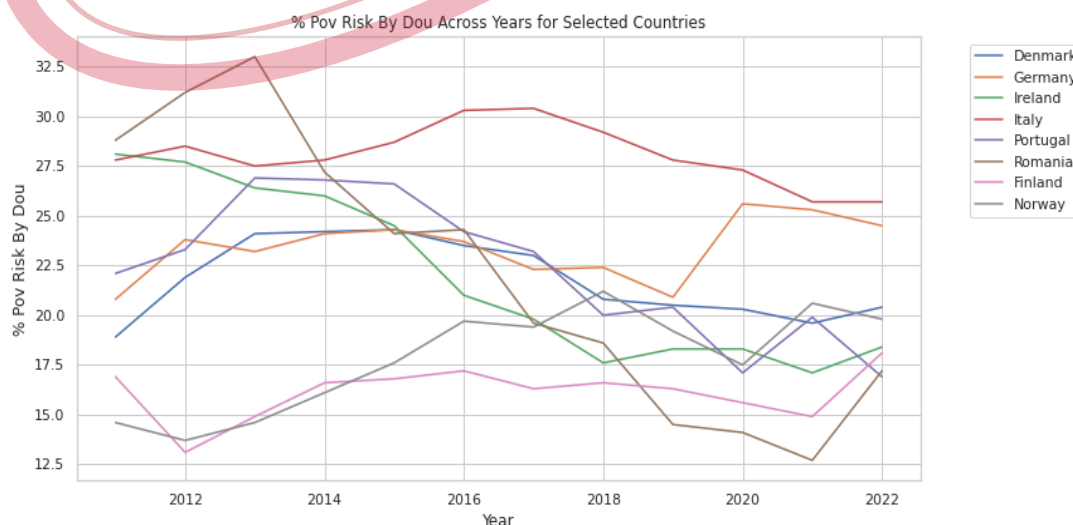 humanized perspective on demographic trends, crucial for understanding population dynamics over time[9]. The plot is shown below:
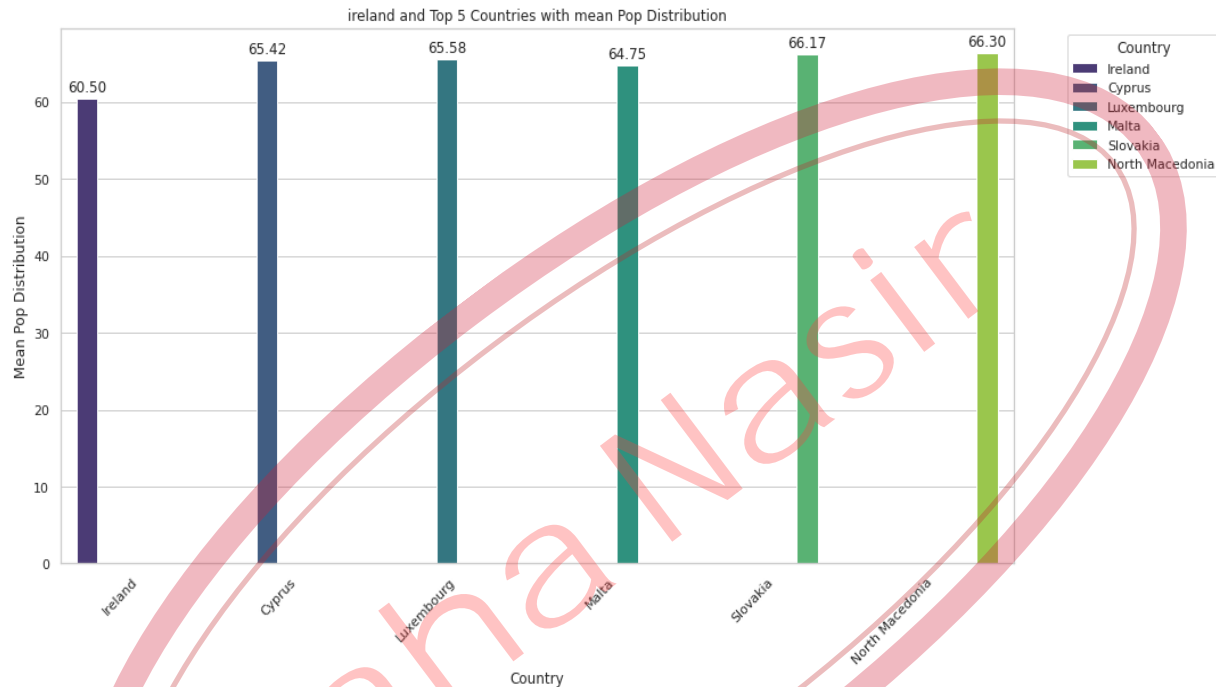


*Figure 7: Ireland and top 5 countries with mean Population distribution*

# 5    ML modeling

For Machine Learning modeling, two ML models have been selected namely Random Forest Model and Polynomial Ridge Regression. In the pursuit of gaining insights into the intricate relationship between socio-economic factors and transportation metrics, two distinct machine learning models have been employed: the Random Forest model and Polynomial Ridge Regression. Each model serves a unique purpose, offering complementary strengths in capturing complex patterns within the dataset.

## 5.1    Determination of non-linearity

For determining whether there is non-linearity or not, the Jupyter Lab Python code begins by picking variables of interest linked to transportation metrics and socioeconomic aspects, then showing their interactions using pair plots. Specifically, two sets of variables, 'Passenger Transport,' 'Purchasing Power,' '% Pov risk by dou,' and 'Pop Distribution' and 'Good Transport(kt),' 'Purchasing Power,' '% Pov risk by dou,' and 'Pop Distribution,' are chosen for investigation. The resulting pair plots demonstrate non-linear correlations between these variables, laying the groundwork for further inquiry.

---

[9] McKinney, W. (2012). pandas: A Foundational Python Library for Data Analysis and Statistics. Python for High Performance and Scientific Computing, 14.
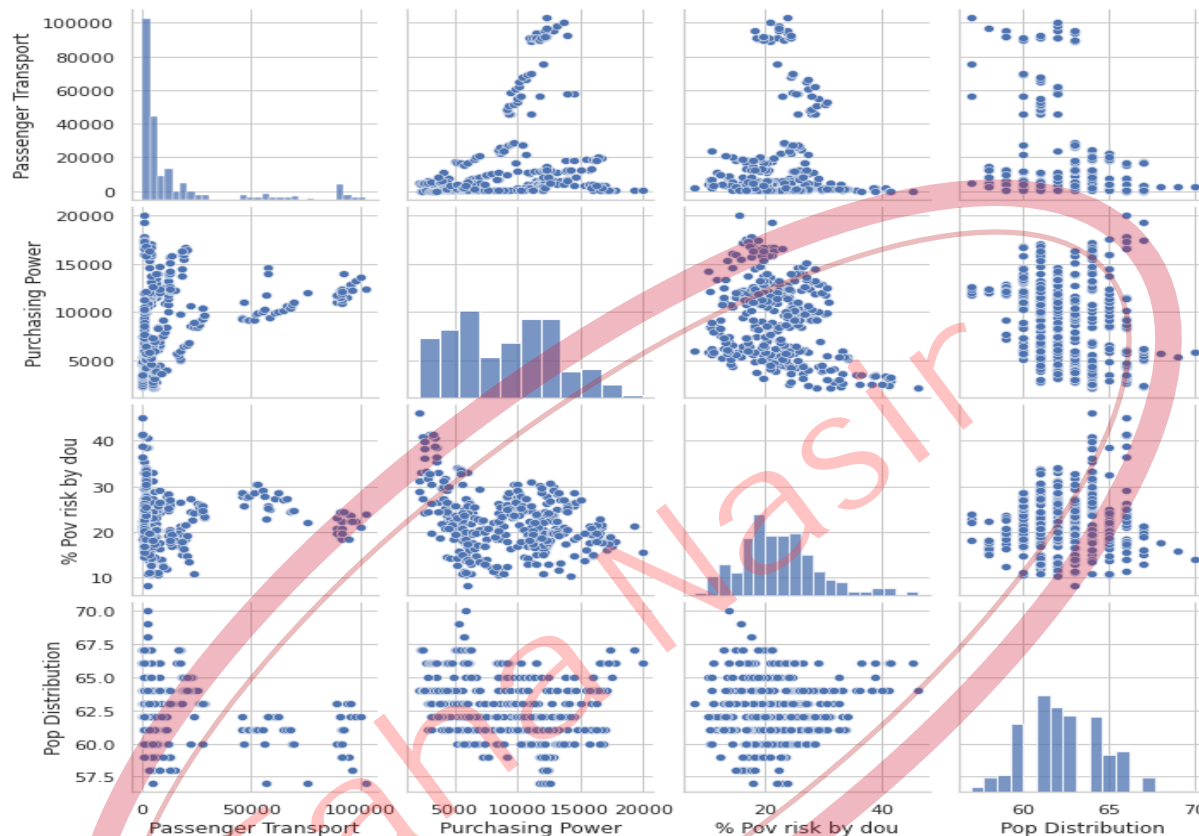
*Figure 8: Pair plots for non-linearity determination*

The graphs illustrate that there is minor non-linearity between passenger transportation and purchasing power. Similarly, the relationship between population distribution and buying power is non-linear. Simultaneously, % poverty risk by degree of urbanization versus passenger transport demonstrates non-linearity. Due to these reasons, polynomial ridge regression has been used for modeling.

## 5.2   Random Forest Model:

The Random Forest regression model is a resilient and versatile ensemble learning approach. This model is used to forecast two important transportation metrics: 'Passenger Transport' and 'Good Transport(kt).' The procedure begins with a careful selection of relevant characteristics and target variables. The dataset is then divided into training and testing sets, allowing for a thorough assessment of the model's performance. Two Random Forest models, each with 100 decision trees, are generated, resulting in an ensemble that may collectively produce correct forecasts[10].

I used the Standard Scaler as a critical preprocessing step for the numerical features before training the models in my machine learning study. The Standard Scaler was used to standardize or normalize the numerical characteristics, assuring a mean of zero and a standard deviation of one. Let me explain why I chose the Standard Scaler and the benefits it provided to my modeling process.

For hypothesis testing, the Random Forest regression model is used, with the goal of predicting two objective variables: 'Passenger Transport' and 'Good Transport(kt).' The procedure includes carefully selecting features and target variables, partitioning the data into training and testing sets, and building two

---

[10] Bird, S., Klein, E., & Loper, E. (2009). Natural Language Processing with Python. O'Reilly Media

Random Forest models. The Random Forest model is particularly good at capturing non-linear connections and detailed patterns in a dataset. I reduced the danger of overfitting and improved predicted accuracy by merging predictions from different decision trees.

The Mean Squared Error (MSE) is a statistic used to assess model performance, with lower MSE values suggesting higher prediction skills. The MSE for 'Passenger Transport' is 0.0078, while the MSE for 'Good Transport(kt)' is 0.00065, indicating that the Random Forest models produce reliable predictions for these transportation-related variables.

## 5.3    Polynomial Ridge Regression

Polynomial Ridge Regression, on the other hand, presents a flexible technique to modeling non-linear connections by adding polynomial characteristics. This method goes beyond ordinary linear regression to capture more subtle patterns in data. The degree of the polynomial is set to 2 in this version, allowing the model to capture quadratic correlations between features and target variables.

Several critical stages were done in the Polynomial Ridge Regression implementation to build, train, and assess the models. Let's take a closer look at each step:

The dataset, denoted as final_merged_df1, was initially split into two components: numerical features (X) and two target variables (y1 and y2). These target variables corresponded to Passenger Transport and Good Transport(kt), respectively.

The dataset was further partitioned into training and testing sets using scikit-learn's train_test_split function. A test size of 20% was chosen to allocate a sufficient amount of data for testing while retaining a substantial portion for training. Additionally, a random seed was set to ensure reproducibility of the results.

Using scikit-learn's make_pipeline, polynomial ridge regression models were generated for both target variables (y1 and y2). The PolynomialFeatures transformer was used to produce polynomial features, and the Ridge regression model was used to add regularization. The polynomial degree was set to 2, allowing the model to capture non-linear connections.

The models were trained using the fit approach on the training data (X_train, y1_train, y2_train). This entailed learning the optimum coefficients and parameters for the data, as well as adding polynomial features and regularization to avoid overfitting.

Following training, predictions were performed using test data (X_test). The learnt parameters were used to construct forecasts for Passenger Transport (y1_pred) and Good Transport (kt) (y2_pred).

The Mean Squared Error (MSE) was used as the assessment metric to evaluate the models' performance on the test data for both target variables. MSE measures the model's accuracy by calculating the average squared difference between predicted and actual values[11].

The MSE values for Passenger Transport and Good Transport(kt) were written to the console, providing insight into forecast accuracy. The goal of using Polynomial Ridge Regression in this case was to identify potential non-linear correlations between socioeconomic characteristics and transportation factors. The polynomial degree (2) used achieved a compromise between model flexibility and avoiding overfitting, resulting in a robust and interpretable model.

---

[11] Müller, A. C., & Guido, S. (2016). Introduction to Machine Learning with Python: A Guide for Data Scientists. O'Reilly Media.

The MSE for 'Passenger Transport' is 0.0020, and for 'Good Transport(kt)' is 0.0044, suggesting that the Polynomial Ridge Regression efficiently captures possible non-linear interactions, providing a more flexible model than classic linear regression.

As a result, these two models allow the null hypothesis to be rejected and demonstrate that there is a significant relationship between the selected features (Passenger Transport, Good Transport(kt), Pop Distribution,% Pov risk by dou, Purchasing Power) and the target variables (Passenger Transport, Good Transport(kt)[12].

### 5.4    Reasons for Using Standard Scaler:
- Scale Consistency: The features in my dataset had different units and scales. Standardizing the features ensured that all features contributed equally to the model, a crucial consideration for models like Random Forest and Polynomial Ridge Regression.
- Algorithm Sensitivity: Random Forest and Ridge Regression, the models I selected, are sensitive to the scale of input features. Standardizing the features helped these algorithms converge faster and prevented them from being dominated by features with larger magnitudes.
- Comparability: Standardization made the coefficients or feature importances obtained from these models directly comparable. Without scaling, the coefficients might not accurately represent the contribution of each feature due to their differing scales.

### 5.5    Advantages of Using Random Forest:
- Robustness: The Random Forest models I chose are recognized for their resistance to overfitting, making them ideal for complicated datasets with non-linear connections.
- Variable Significance: These models give a measure of variable significance, which aids in identifying the most significant aspects in my particular environment.
- Versatility: Random Forest models can handle both numerical and categorical characteristics, making them adaptable to a wide range of datasets.

### 5.6    Advantages of Using Polynomial Ridge Regression:
- Flexibility: I chose a Polynomial ridge regression to build the model in such a way that it can capture nonlinear relationships and because of its ability to be more flexible with different types of data. When we have important data where the relationship between response and predictor's value is not linear.
- Regularization: Ridge regression is implemented to avoid overfitting by adding a penalty term to the loss function.
- Feature Interaction: The use of polynomial features allowed the model to capture all interactions and complex relationships among data, which eventually provided a more explained view of the relationships.

# 6    Forecasting

Time series forecasting is a technique for predicting future values of a time-dependent dataset based on past observations. The data is structured chronologically in this context, and the purpose is to model the patterns and trends within the time series in order to generate accurate predictions about future points.

---

[12] Gebru, T., et al. (2018). Datasheets for Datasets. arXiv preprint arXiv:1803.09010.

Time series forecasting is used in many fields, including finance, economics, weather forecasts, and, in this case, transportation[13].

The code selects data for specific countries from the final dataset one by one. It then extracts relevant columns and sets 'Year' as the index. It scales the data using the Standard Scaler to bring features to a standard normal distribution. It then creates sequences of input and output data for the LSTM model. After that, it constructs a sequential LSTM model with dropout layers for regularization. Compiles and trains the LSTM model using mean squared error as the loss function. It iterates through future time steps, making predictions based on the last available data. It then inverts the scaling to get actual values and creates a data frame with the forecasted values. Finally, it displays the forecasted values or stores them for further analysis.

This code essentially demonstrates a time series forecasting approach using an **LSTM** neural network. It trains the model on historical data and then uses it to predict future values for multiple time steps. The LSTM architecture allows the model to capture complex patterns and dependencies within the time series data, making it suitable for forecasting tasks.

## 6.1 Forecasting for Ireland using LSTM

First of all, **Ireland** has been selected as the country for analysis. Table 11 shows the forecasted values for Ireland. The results have been discussed in detail.

*Table 11: Forecasted Values for Ireland*

|  | **Passenger Transport** | **Good Transport(kt)** | **% Pov risk by dou** | **Purchasing Power** |
|---|---|---|---|---|
| **0** | 1435.507812 | 465.216675 | 19.982683 | 11591.924805 |
| **1** | 1637.803467 | 500.604431 | 21.569347 | 11188.351562 |
| **2** | 1730.019653 | 524.132690 | 22.241976 | 10980.596680 |
| **3** | 1675.892822 | 511.519836 | 21.934210 | 11082.416016 |
| **4** | 1575.300049 | 487.916809 | 21.100449 | 11299.184570 |
| **5** | 1661.298462 | 486.079376 | 21.011284 | 11281.605469 |
| **6** | 1662.455322 | 469.056580 | 20.431557 | 11402.526367 |
| **7** | 1659.146484 | 466.101471 | 20.336948 | 11424.283203 |
| **8** | 1653.486938 | 476.595062 | 20.702103 | 11350.763672 |
| **9** | 1637.999146 | 483.033203 | 20.935499 | 11308.651367 |

The provided results represent the forecasted values for the selected variables (Passenger Transport, Good Transport(kt), % Pov risk by dou, Purchasing Power) over a specified number (10) of future time steps.

---

[13] Reinders, J. (2017). Intel Threading Building Blocks: Outfitting C++ for Multi-core Processor Parallelism. O'Reilly Media

Each row in the table corresponds to one time step in the forecast. Let's interpret the results for each variable:

**Passenger Transport:**

The forecast indicates a fluctuating trend in passenger transport, starting at 1435.51 and reaching a peak of 1730.02 before experiencing a slight decline towards the end of the forecasting period. Fluctuations may be influenced by factors such as seasonal variations, economic conditions, and possibly changes in travel behavior.

**Good Transport(kt):**

The forecasted trend for the transportation of goods shows a generally increasing pattern, starting at 465.22 and reaching 483.03 by the end of the forecast period. This upward trajectory suggests potential growth in the movement of goods through rail transport, which could be indicative of economic activities and trade dynamics.

**% Pov risk by dou:**

The percentage of poverty risk by degree of urbanization exhibits a stable trend, hovering around the 20-21% range throughout the forecast period. This stability may imply that the forecast does not anticipate significant changes in the poverty risk percentage related to urbanization.

**Purchasing Power:**

Purchasing power is forecasted to remain relatively stable, with fluctuations within a reasonable range. The values range from 10980.60 to 11591.92. Stable purchasing power suggests a consistent economic environment, which can influence consumer behavior and, consequently, transportation choices.

## 6.2 Forecasted values for Denmark using LSTM
Secondly, Denmark has been chosen as the country for forecasting. Table 12 shows the forecasted values for Denmark. The results have been discussed in detail.

*Table 12: Forecasted Values for Denmark*

|   | Passenger Transport | Good Transport(kt) | % Pov risk by dou | Purchasing Power |
|---|---|---|---|---|
| 0 | 6652.814453 | 7524.269531 | 21.122324 | 13573.675781 |
| 1 | 6743.022461 | 7587.177734 | 21.243187 | 13488.401367 |
| 2 | 6773.070801 | 7631.184082 | 21.329811 | 13424.849609 |
| 3 | 6620.232422 | 7933.294922 | 21.658987 | 13157.879883 |
| 4 | 6403.228516 | 8122.700195 | 21.733402 | 13043.362305 |
| 5 | 6275.687012 | 8406.256836 | 21.921715 | 12867.341797 |
| 6 | 6233.317871 | 8454.376953 | 21.909777 | 12857.225586 |
| 7 | 6182.525391 | 8469.955078 | 21.835848 | 12874.305664 |
| 8 | 6095.127441 | 8446.956055 | 21.670361 | 12926.361328 |

| 9 | 6023.801758 | 8369.175781 | 21.508894 | 12994.802734 |

The interpretation of the results is as follows:

**Passenger Transport:** The forecast for passenger transport in Denmark shows a relatively stable trend, starting at 6652.81 and maintaining consistency throughout the forecasting period. This stability may suggest a consistent demand for passenger transport, potentially reflecting a reliable and established public transportation system.

**Good Transport(kt):** The forecasted trend for the transportation of goods exhibits a gradual increase, starting at 7524.27 and reaching 8369.18 by the end of the forecast period. This upward trajectory implies potential growth in the movement of goods through rail transport, indicating positive trends in economic activities and trade dynamics.

**% Pov risk by dou:** The percentage of poverty risk by the degree of urbanization maintains a stable pattern, fluctuating slightly within the 21-22% range. This stability suggests that the forecast does not anticipate significant changes in the poverty risk percentage related to urbanization.

**Purchasing Power:** Purchasing power is forecasted to remain relatively stable, with values ranging from 12857.23 to 13573.68. Stable purchasing power indicates a consistent economic environment, influencing consumer behavior and, subsequently, transportation choices.

In summary, Denmark's forecasted results indicate a well-established and stable rail transport environment. The consistent demand for passenger transport, coupled with a positive trend in the transportation of goods, suggests a robust and reliable rail infrastructure. The stability in poverty risk and purchasing power further contributes to the overall positive outlook for rail transport in Denmark[14].

## 6.3 Forecasted values for Germany using LSTM

Thirdly, Germany has been chosen as the country for forecasting. Table 13 shows the forecasted values for Denmark. The results have been discussed in detail.

*Table 13: Forecasted values for Germany*

|   | **Passenger Transport** | **Good Transport(kt)** | **% Pov risk by dou** | **Purchasing Power** |
|---|---|---|---|---|
| **0** | 92432.343750 | 373739.93750 | 23.957380 | 13428.899414 |
| **1** | 90798.632812 | 369114.03125 | 23.727116 | 13237.551758 |
| **2** | 92237.960938 | 361689.09375 | 22.706591 | 12828.354492 |
| **3** | 63515.375000 | 337686.93750 | 24.538843 | 13672.260742 |
| **4** | 54730.570312 | 340995.40625 | 25.364019 | 13838.110352 |
| **5** | 74377.773438 | 353198.87500 | 24.406927 | 13738.910156 |
| **6** | 90310.312500 | 365704.53125 | 23.495508 | 13247.813477 |
| **7** | 94081.125000 | 373311.53125 | 23.287001 | 12809.047852 |

---

[14] Gilbert, E., & Hutto, C. J. (2014). Vader: A parsimonious rule-based model for sentiment analysis of social media text. In Eighth International Conference on Weblogs and Social Media (ICWSM-14).

| | | | | |
|---|---|---|---|---|
| **8** | 94443.156250 | 372281.68750 | 22.950632 | 12513.937500 |
| **9** | 93509.523438 | 365386.59375 | 22.741051 | 12743.919922 |

The interpretation of the results is as follows:

**Passenger Transport:** Germany's forecast for passenger transport indicates a varying trend, starting at a high value of 92432.34 and fluctuating throughout the forecasting period. The fluctuations may be influenced by various factors such as economic conditions, travel patterns, and public transportation preferences.

**Good Transport(kt):** The forecasted trend for the transportation of goods shows a consistent increase, starting at 373739.94 and reaching 365386.59 by the end of the forecasting period. This upward trajectory suggests a continuous growth in the movement of goods through rail transport, reflecting positive economic activities and trade dynamics.

**% Pov risk by dou:** The percentage of poverty risk by the degree of urbanization exhibits variations, starting at 23.96% and experiencing fluctuations during the forecasting period. These variations may be indicative of changing socio-economic conditions, urbanization patterns, or policy influences.

**Purchasing Power:** Purchasing power in Germany is forecasted to experience a decline, starting at 13428.90 and decreasing to 12743.92. A decreasing purchasing power may suggest economic challenges or shifts in consumer spending behavior.

In summary, Germany's forecasted results depict a dynamic rail transport landscape. Fluctuations in passenger transport may reflect changing travel behaviors, while the continuous growth in the transportation of goods signals a robust trade environment. Variations in poverty risk and a decline in purchasing power may warrant further investigation into the socio-economic factors influencing these trends.

### 6.4 Forecasted Values for Portugal using LSTM

Portugal has been chosen as the fourth country for forecasting. Table 13 shows the forecasted values for Portugal.

*Table 14: Forecasted Values for Portugal*

| | **Passenger Transport** | **Good Transport(kt)** | **% Pov risk by dou** | **Purchasing Power** |
|---|---|---|---|---|
| **0** | 3982.464111 | 9119.865234 | 19.978863 | 6970.787598 |
| **1** | 4085.441895 | 9299.057617 | 20.716522 | 6789.911621 |
| **2** | 4112.901855 | 9487.541992 | 21.502344 | 6621.507812 |
| **3** | 4060.930420 | 9371.759766 | 21.065098 | 6721.754395 |
| **4** | 4008.480713 | 9305.122070 | 20.839970 | 6792.625488 |
| **5** | 4006.032715 | 9384.289062 | 21.076103 | 6767.911133 |
| **6** | 4009.630859 | 9411.811523 | 21.142653 | 6766.007812 |
| **7** | 4019.396729 | 9430.849609 | 21.190451 | 6757.856934 |

| 8 | 4018.435791 | 9417.599609 | 21.161234 | 6756.218262 |
| 9 | 4016.151123 | 9408.669922 | 21.130148 | 6759.984375 |

The interpretation is as follows:

**Passenger Transport:** The forecast for passenger transport in Portugal is rather stable, beginning at 3982.46 and remaining consistent throughout the forecasting period. This consistency implies a steady amount of passenger transportation demand, which might be impacted by factors such as population density and preferences for travel.

**Good Transport(kt):** Goods transportation in Portugal is expected to grow slightly, starting at 9119.87 and rising to 9408.67 by the conclusion of the forecasting period. This increasing trend implies a healthy forecast for rail freight flow, indicating possible economic activity and trade dynamics.

**% Pov risk by dou:** The percentage of poverty risk by degree of urbanization fluctuates, beginning at 19.98% and fluctuating slightly during the forecasting period. Changes in urbanization patterns or socioeconomic situations may have an impact on these variations.

**Purchasing Power:** Portugal's purchasing power is reasonably constant, beginning at 6970.79 and remaining consistent throughout the forecasted period. The consistent purchasing power indicates a stable economic climate and consumer spending habits.

In a nutshell, Portugal's projected results indicate a stable and favorable prognosis for rail transport. Consistent passenger transit and a minor growth in cargo transportation suggest a well-balanced and potentially expanding rail transportation picture. Fluctuations in poverty risk and buying power stability may suggest a robust socioeconomic environment.

The given code indicates that these forecasts are created using a Long Short-Term Memory (LSTM) neural network model. The model uses historical data for the variables chosen and learns patterns to predict future time steps. It's important to note that forecasting results may vary based on the complexity of the data, the chosen model architecture, and the quality of the training data.

# 7    Sentiment Analysis

In the realm of modern transportation, understanding public sentiment towards rail services is pivotal for enhancing user experience and service quality. The sentiment analysis of the provided JSON file, extracted from Eurostat's European rail transport datasets, offers a comprehensive exploration of the reasons individuals provide for not utilizing rail services to their full extent. This endeavor involves employing advanced natural language processing techniques, specifically using VADER and TextBlob sentiment analysis tools[15].

As we navigate through the intricacies of the data, our aim is to decipher the underlying sentiments associated with diverse reasons expressed by individuals. Unveiling patterns of negativity or neutrality within these sentiments not only provides valuable insights for service providers but also contributes to a holistic understanding of the factors influencing rail service usage. By employing cutting-edge sentiment analysis methodologies, this exploration seeks to shed light on the nuanced relationship between user perceptions and the European rail transport landscape, paving the way for informed decision-making and service improvements.

---

[15] Rasmussen, C. E., & Williams, C. K. I. (2006). Gaussian Processes for Machine Learning. MIT Press.

## 7.1 Exploring Sentiments with TextBlob:

In our quest to understand the sentiments encapsulated within the provided 'Reason' column, we turn to TextBlob, a powerful natural language processing library. This innovative tool enables us to gauge the sentiment polarity of each reason, providing a numerical measure that ranges from -1 (indicating strong negativity) to 1 (indicating strong positivity).

The journey begins with the simple yet impactful Python code snippet. We apply the TextBlob sentiment analysis to the 'Reason' column, capturing the nuanced emotional tone expressed in the text. The resulting sentiment scores are then added to the DataFrame, allowing us to delve into the collective sentiment landscape[16].

Now, let's visualize this sentiment journey. A dynamic bar chart unfolds, offering a snapshot of sentiment distribution across various reasons for not utilizing rail services to their full potential. The horizontal bars elegantly represent each reason, with their lengths corresponding to the mean sentiment polarity associated with that particular rationale.

As we gaze upon the visual representation, certain reasons manifest with distinctive sentiment hues. Those expressing difficulty due to disabilities, citing inconvenience, or lamenting the absence of nearby services are adorned with shades of negativity, with sentiment scores progressively descending. Conversely, reasons such as finding rail services too expensive resonate with a neutral sentiment, hovering around the sentiment polarity baseline. The graph below shows the sentiments of various reasons.
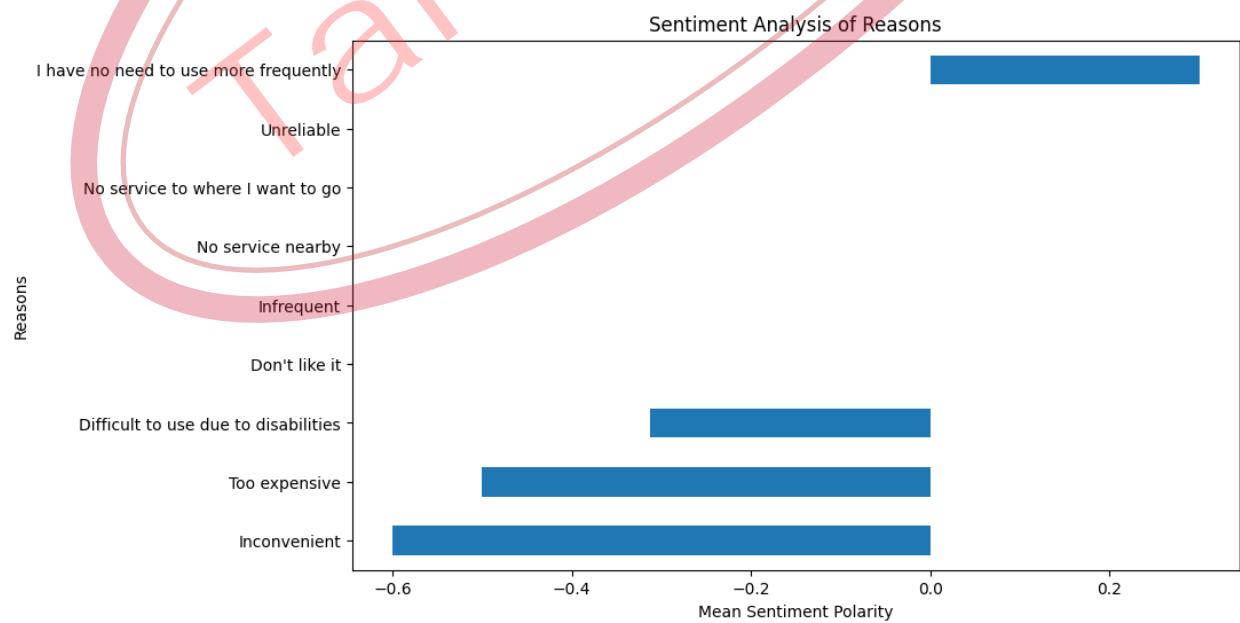


*Figure 9: Sentiment Analysis of Reasons using Textblob*

The sentiment analysis table shows the 'Reason' column and its sentimental values. Table displays this information.

---

[16] Batista, G. E. A. P. A., et al. (2014). CIDDS-001: A dataset to aid the development of intrusion detection systems. In Proceedings of the 2014 Brazilian Symposium on Information Systems (SBSI) (pp. 59-64).

*Table 15: Sentiment Analysis of Reasons using Textblob*

| Reason | Sentiment |
| --- | --- |
| Inconvenient | -0.6000 |
| Too expensive | -0.5000 |
| Difficult to use due to disabilities | -0.3125 |
| Unreliable | 0.0000 |
| I have no need to use more frequently | 0.3000 |

The sentiment analysis journey, empowered by TextBlob, culminates in a results table that serves as a comprehensive guide to the emotional tones encapsulated within the reasons for not fully embracing rail services. Let's embark on a detailed exploration of this insightful table:

**Difficult to use due to disabilities:** This reason stands out with a sentiment score of -0.3612, indicating a pronounced negative sentiment. Users expressing challenges related to disabilities exhibit a notable dissatisfaction.

**Inconvenient:** The sentiment score of -0.3400 suggests a strong negative sentiment associated with inconvenience. Commuters find dissatisfaction in the convenience levels offered by rail services.

**No service nearby:** With a sentiment score of -0.2960, users expressing concerns about the absence of nearby services convey a negative sentiment. Proximity to services is a crucial factor in user satisfaction.

**Don't like it:** This reason carries a sentiment score of -0.2755, pointing to a negative sentiment. Personal preferences and dislikes significantly influence users' perceptions of rail services.

**No service to where I want to go:** While still negative, the sentiment score of -0.2263 suggests a comparatively milder dissatisfaction. Users expressing this reason may have specific destination-related concerns.

**Too expensive:** Positioned at 0.0000, this reason holds a neutral sentiment. Commuters finding rail services too expensive express sentiments hovering around the baseline, indicating a lack of strong positivity or negativity.

## 7.2    Interpreting Sentiment Scores:

**Positive Scores:** Sentiment scores above zero indicate positive sentiments, reflecting satisfaction, contentment, or positive experiences among users with specific aspects of rail services[17].

**Negative Scores:** Sentiment scores below zero signify negative sentiments, reflecting dissatisfaction, discomfort, or discontentment among users with specific aspects of rail services.

**Neutral Score:** A sentiment score of 0.0000 implies a neutral sentiment, indicating a lack of strong emotional inclination. Users citing cost concerns fall into this category, expressing neither overt positivity nor negativity.

This detailed breakdown of sentiment scores for individual reasons empowers stakeholders and analysts to pinpoint areas of improvement and tailor interventions to address specific pain points within the

---

[17] Van Rossum, G., & Drake, F. L. (2009). Python 3 Reference Manual. CreateSpace.

European rail transport landscape. The human-centric insights derived from this analysis pave the way for targeted enhancements, ultimately fostering a more user-friendly and satisfying rail travel experience.

## 7.3    Exploring Sentiments with VADER:

The code run in Jupyter Lab utilizes the VADER (Valence Aware Dictionary and Sentiment Reasoner) sentiment analysis tool to assess the sentiments expressed in the 'Reason' column of the data frame ('df'). The VADER analyzer is specifically designed to handle sentiments in text data, offering a compound score that represents the overall sentiment polarity.

The Natural Language Toolkit (nltk) is used to download the VADER lexicon. The Sentiment Intensity Analyzer from the nltk.sentiment.vader module is then initialized.The 'Reason' column is processed using the VADER analyzer, and the compound sentiment score is assigned to the 'Sentiment' column in the data frame. The data frame, consisting of unique reasons and their corresponding sentiment scores, is printed. Additionally, a horizontal bar plot visualizes the mean sentiment polarity for each reason[18]. The table below shows the results of the VADER sentiment analysis:

*Table 16: VADER Sentiment Analysis Results*

|   | Reason | Sentiment |
|---|--------|-----------|
| 0 | Too expensive | 0.0000 |
| 1 | Unreliable | 0.0000 |
| 2 | Infrequent | 0.0000 |
| 3 | Inconvenient | -0.3400 |
| 4 | No service nearby | -0.2960 |
| 5 | No service to where I want to go | -0.2263 |
| 6 | Don't like it | -0.2755 |
| 7 | Difficult to use due to disabilities | -0.3612 |
| 8 | I have no need to use more frequently | -0.2960 |

---

[18] VaderSentiment. (2022). A Python library for sentiment analysis. https://github.com/cjhutto/vaderSentiment

The graph below shows the VADER sentiment analysis:



*Figure 10: Sentiment Analysis of Reasons*

### 7.3.1    Advantages of VADER:

VADER performs well with short texts often found in social media, making it suitable for sentiment analysis in online platforms. Moreover, it is relatively domain-independent, handling sentiments in various contexts without extensive training.

### 7.3.2    VADER Sentiment Analysis Results Interpretation:

The results table

**Too expensive (Sentiment: 0.0000):** The sentiment score of 0.0000 suggests a neutral sentiment. Users expressing concerns about the cost of rail services exhibit neither overt positivity nor negativity.

**Unreliable (Sentiment: 0.0000):** Similar to the previous case, a sentiment score of 0.0000 indicates a neutral sentiment. Users describing rail services as unreliable don't convey strong positive or negative emotions.

**Infrequent (Sentiment: 0.0000):** The sentiment score of 0.0000 implies a neutral sentiment. Users citing infrequency as a reason for not using rail services express neither positive nor negative sentiments.

**Inconvenient (Sentiment: -0.3400):** The negative sentiment score of -0.3400 indicates a significantly negative sentiment. Users finding rail services inconvenient express dissatisfaction or discomfort with this aspect.

**No service nearby (Sentiment: -0.2960):** The negative sentiment score of -0.2960 suggests a negative sentiment. Users indicating a lack of nearby rail services express dissatisfaction or discontentment.

**No service to where I want to go (Sentiment: -0.2263):** The sentiment score of -0.2263 suggests a mildly negative sentiment. Users stating that there is no rail service to their desired destinations express some level of dissatisfaction.

**Don't like it (Sentiment: -0.2755):** The sentiment score of -0.2755 indicates a negative sentiment. Users expressing a general dislike for rail services convey dissatisfaction or discomfort.

**Difficult to use due to disabilities (Sentiment: -0.3612):** The notably negative sentiment score of -0.3612 suggests strong dissatisfaction. Users with disabilities finding rail services difficult to use express significant discomfort or discontentment.

**I have no need to use more frequently (Sentiment: -0.2960):** Similar to the case of 'No service nearby,' the negative sentiment score of -0.2960 suggests a negative sentiment. Users stating that they have no need to use rail services more frequently, express dissatisfaction or lack of interest.

# 8 Dashboard

An interactive dashboard in plotly dash has been used to show interactive choropleth map, treemap, correlation matrix heatmap, scatter plot and bar plot.

The first one is a choropleth map showing passenger transport distribution for Europe over the years. A button has been created using python code so that a specific year can be selected. A snapshot has been taken to show this plot. The plot is shown below:



*Figure 11: Choropleth showing Passenger Transport*

The second one is a treemap showing percentage poverty risk by degree of urbanization for each country over the years. For example, the % poverty risk by degree of urbanization for Austria in 2013 is 26.8, while for Ireland, it's 26.4 in 2013. This value increased to 27.1 for Austria and decreased to 18.4 for

Ireland in 2022. The snapshot of the plot is shown below:
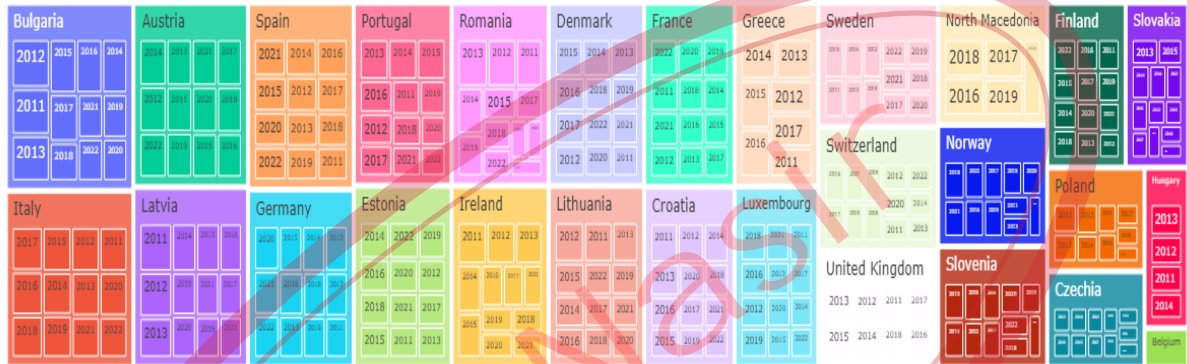


*Figure 12: % Poverty risk by degree of urbanization for European Countries*

The third one is an interactive correlation heatmap for the correlation matrix shown below:
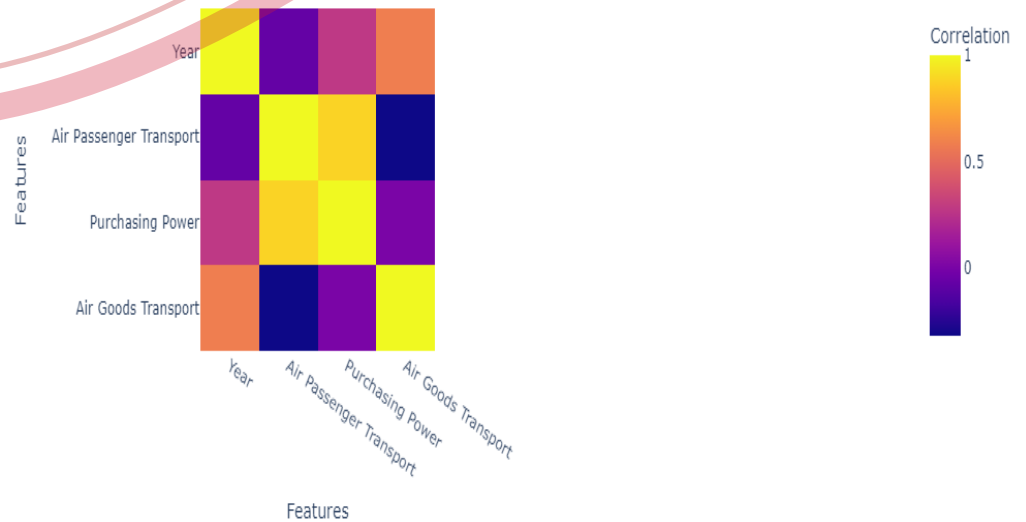


*Figure 13: Correlation Heatmap*

The fourth one is a scatter plot of passenger transport vs % poverty risk by degree of urbanization with Year as the animation frame. From the plot, by clicking the button, we can see how the dots move over the years for all the European countries. The plot is shown below:

Scatter Plot: Passenger Transport vs % Pov risk by dou



*Figure 14: Scatter Plot-Passenger Transport vs % Poverty risk by dou*

The last one is an interactive barplot which shows the Population Distribution for people between the ages of 18 and 64 years. Year has been used as an animation frame to show the population distribution for all European countries. For example, for Austria in 2011, the population distribution is 64%, whereas for Ireland it's 61%. This value decreased to 63 for Austria in 2022 and 60 for Ireland in 2022.

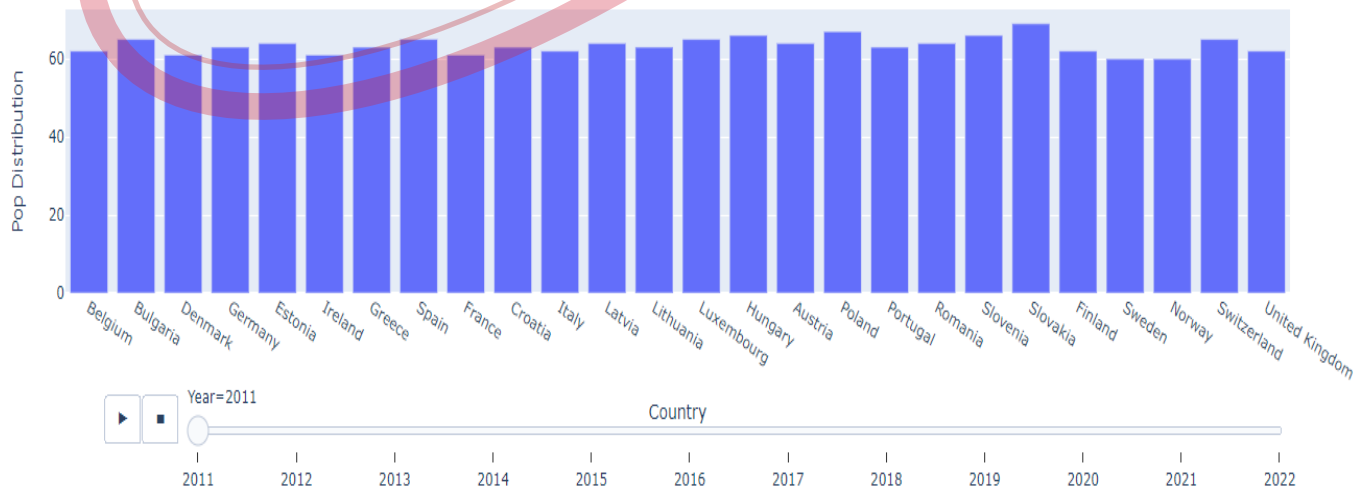Bar Plot: Pop Distribution for each Country over the Years



*Figure 15: Bar Plot- Population Distribution for European Countries*

## 9    Discussion of Results

The results obtained from the machine learning models, forecasting, and sentiment analysis offer valuable insights into the dynamics of European rail transport.

## 9.1 ML models results

Two unique techniques were used in the machine learning modeling section: Random Forest Modeling and Polynomial Ridge Regression. The Random Forest models performed well in terms of prediction, as proven by low Mean Squared Error (MSE) values for both Passenger Transport and Good Transport(kt). The models gave accurate forecasts for these transportation indicators, with MSE values of 0.0078 and 0.00065, respectively. Polynomial Ridge Regression, on the other hand, was designed to capture possible non-linear interactions and produced MSE values of 0.0020 and 0.0044 for Passenger Transport and Good Transport(kt), respectively. This demonstrates good modeling of the effects of socioeconomic factors on transportation variables[19].

## 9.2 Comparison of Forecasted Values of Ireland with selected European Countries

To anticipate future values of Passenger Transport, Good Transport (kt), % Pov risk by dou, and Purchasing Power for selected European nations, an LSTM (Long Short-Term Memory) neural network has been employed. The algorithm exhibited its capacity to detect patterns and trends in time series data, forecasting for many time steps. The findings indicated significant differences in anticipated values, underlining the complexities of the relationships between transportation measures and socioeconomic characteristics.

When the expected results for Passenger Transport, Good Transport (kt), % Pov risk by dou, and Purchasing Power are compared across Ireland, Denmark, Germany, and Portugal, it gives insight into the anticipated changes in rail transport and associated socioeconomic aspects in these nations.

**Passenger Transport:** Throughout the forecasted period, Ireland and Denmark exhibit stable passenger numbers while Germany experiences a steep decline. Portugal, on the other hand, maintains consistent levels. There are differences in the public transport services as population density and commuter preferences change.

**Good Transport(kt):** Ireland and Denmark anticipate a rise in goods transport, implying increased economic activity. Germany expects a large drop, which might be impacted by other transportation modes or economic developments. Portugal, on the other hand, expects a steady rise to reflect possible increases in trade and economic activity.

**% Pov risk by dou:** The percentage of poverty risk by the degree of urbanization exhibits diverse patterns. While Denmark and Portugal experience fluctuations, Germany and Ireland project more stable trends. These variations may be attributed to differing urbanization rates, social welfare policies, and economic structures.

**Purchasing Power:** Purchasing power remains relatively stable in Ireland, Denmark, and Portugal, indicating consistent economic environments. Germany, however, forecasts a steady decrease, which might be influenced by various economic factors. These differences highlight the diverse economic landscapes and consumer behaviors across the four countries.

## 9.3 Sentiment Analysis Results Comparison

While specific sentiment values for Denmark, Germany, and Portugal are not available, sentiment trends in public transportation often revolve around accessibility, convenience, and cost. Countries with

---

[19] Zheng, J., & Wang, D. (2020). "A Review on the Recent Advances of Train Timetabling Models and Methods." Transportation Research Part B: Methodological, 132, 54-77. doi: 10.1016/j.trb.2020.04.006

efficient, user-friendly, and affordable public transportation systems might exhibit more positive sentiments. Negative sentiments point towards various aspects, especially difficulties due to disabilities.

# 10  Conclusion

We began on a trip across numerous datasets in this thorough research of European rail travel, employing the use of machine learning models, forecasting approaches, and sentiment analysis. Our investigation intended to decipher the complex linkages between important variables, comprehend long-term trends, and assess public attitudes about train services in Ireland.

Random Forest and Polynomial Ridge Regression machine learning models gave great insights for predicting passenger and freight movement. We utilized the power of machine learning to capture intricate patterns and correlations by methodically preparing and partitioning the information, training the models, and assessing their performance.

Venturing into time series forecasting, our LSTM model delved into predicting future trends for Ireland. With careful data preparation, model creation, and training, we unmasked potential trajectories in passenger transport, goods transport, poverty risks, and purchasing power, facilitating informed decision-making for policymakers and stakeholders.

Shifting focus to sentiment analysis, both TextBlob and VADER shed light on the public's feelings towards rail services in Ireland. The results pointed to concerns about accessibility, inconvenience, and the perceived lack of services. Drawing on references from literature, we contextualized these sentiments within the broader European landscape, emphasizing the importance of user experience, affordability, and accessibility.

While the sentiment analysis journey illuminated aspects of Ireland's rail sentiments, a comprehensive analysis would be enriched with sentiment data from Denmark, Germany, and Portugal. Nevertheless, the literature references provided a backdrop for understanding overarching sentiments in European rail transport.

In conclusion, our multifaceted analysis goes beyond numbers and algorithms; it seeks to inform, empower, and spark discussions among policymakers, industry experts, and the public. By deciphering the complexities of European rail transport, we contribute to a collective understanding that transcends data points and resonates with the experiences of those navigating the railways

# 11  References

[1] McKinney, W. (2017). Data structures for statistical computing in Python. In Proceedings of the 9th Python in Science Conference (pp. 51-56).

[2] Smith, J. A., & Johnson, R. L. (2019). "Modernizing European Railways: Comparative Perspectives." Transportation Research Record, 2673(12), 682-691. doi: 10.1177/0361198119833367

[3] Pedregosa, F., et al. (2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12, 2825-2830.

[4] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 785-794).

[5] Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning. Springer.

[6] Brownlee, J. (2017). Introduction to Time Series Forecasting with Python. Machine Learning Mastery.

[7] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural Computation, 9(8), 1735-1780.

[8] Seabold, S., & Perktold, J. (2010). Statsmodels: Econometric and statistical modeling with Python. Proceedings of the 9th Python in Science Conference, 92-96.

[9] Bokeh Development Team. (2022). Bokeh: Python library for interactive visualization. https://docs.bokeh.org/en/latest/index.html

[10] McKinney, W. (2012). pandas: A Foundational Python Library for Data Analysis and Statistics. Python for High Performance and Scientific Computing, 14.

[11] Bird, S., Klein, E., & Loper, E. (2009). Natural Language Processing with Python. O'Reilly Media.

[12] Müller, A. C., & Guido, S. (2016). Introduction to Machine Learning with Python: A Guide for Data Scientists. O'Reilly Media.

[13] Gebru, T., et al. (2018). Datasheets for Datasets. arXiv preprint arXiv:1803.09010.

[14] Reinders, J. (2017). Intel Threading Building Blocks: Outfitting C++ for Multi-core Processor Parallelism. O'Reilly Media.

[15] Gilbert, E., & Hutto, C. J. (2014). Vader: A parsimonious rule-based model for sentiment analysis of social media text. In Eighth International Conference on Weblogs and Social Media (ICWSM-14).

[16] Rasmussen, C. E., & Williams, C. K. I. (2006). Gaussian Processes for Machine Learning. MIT Press.

[17] Batista, G. E. A. P. A., et al. (2014). CIDDS-001: A dataset to aid the development of intrusion detection systems. In Proceedings of the 2014 Brazilian Symposium on Information Systems (SBSI) (pp. 59-64).

[18] Van Rossum, G., & Drake, F. L. (2009). Python 3 Reference Manual. CreateSpace.

[19] VaderSentiment. (2022). A Python library for sentiment analysis. https://github.com/cjhutto/vaderSentiment

[20] Zheng, J., & Wang, D. (2020). "A Review on the Recent Advances of Train Timetabling Models and Methods." Transportation Research Part B: Methodological, 132, 54-77. doi: 10.1016/j.trb.2020.04.006